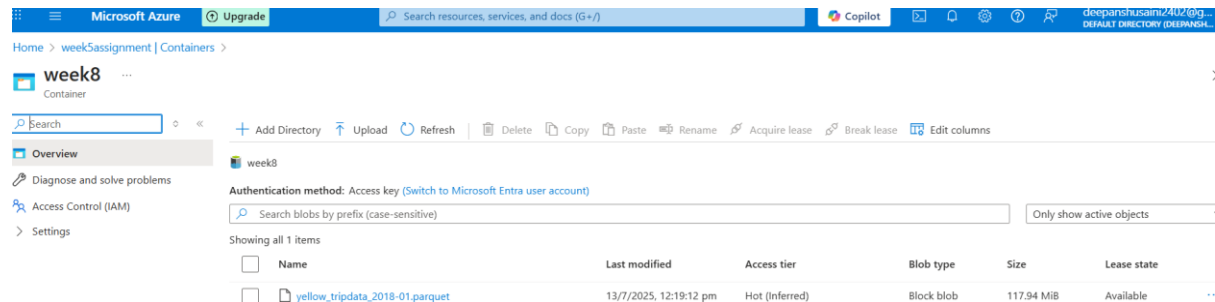# CELEBAL TECHNOLOGIES PVT. LTD.
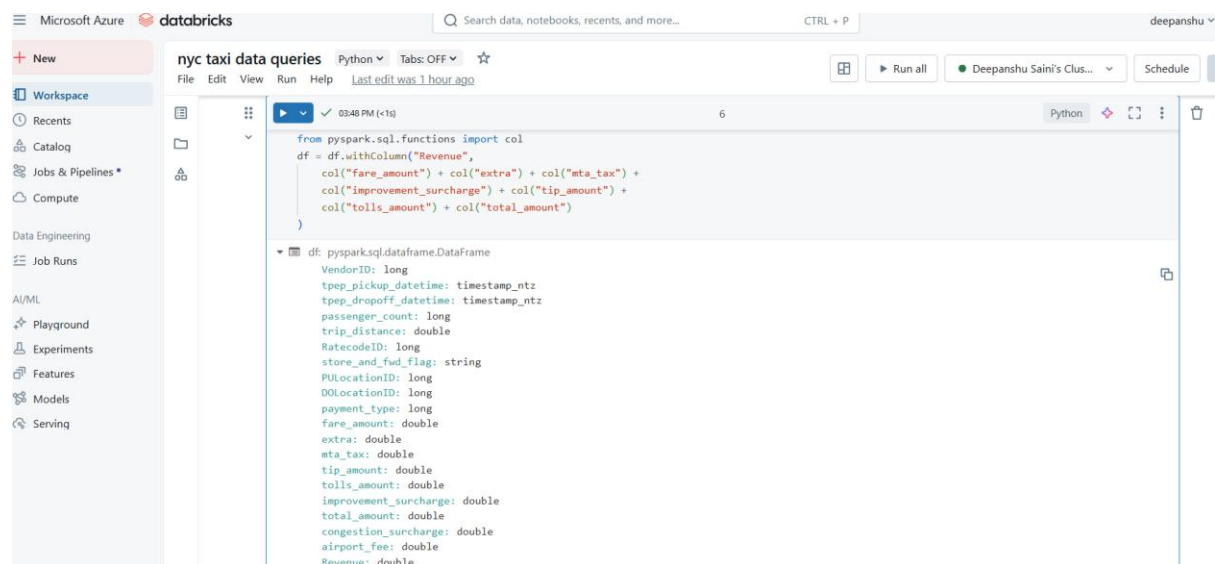
## ASSIGNMENT

### WEEK 8

Load NYC taxi data to DataLake/Blob_Storage/DataBricks



Query 1. - Add a column named as ""Revenue"" into dataframe which is the sum of the below columns:  'Fare_amount', 'Extra', 'MTA_tax', 'Improvement_surcharge', 'Tip_amount', 'Tolls_amount', 'Total_amount'



Query 2. - Increasing count of total passengers in New York City by area

Query 3. - Realtime Average fare/total earning amount earned by 2 vendors



Query 4. - Moving Count of payments made by each payment mode



Query 5. - Highest two gaining vendor's on a particular date with no of passenger and total distance by cab

Query 6. - Most no of passenger between a route of two location.



Query 7. - Get top pickup locations with most passengers in last 5/10 seconds.



Load any dataset into DBFS

## Flatten JSON fields



## Write flattened file as external parquet table